

OCTO：分布式服务通信框架及服务治理系统

张熙 基础架构中心 2016.6

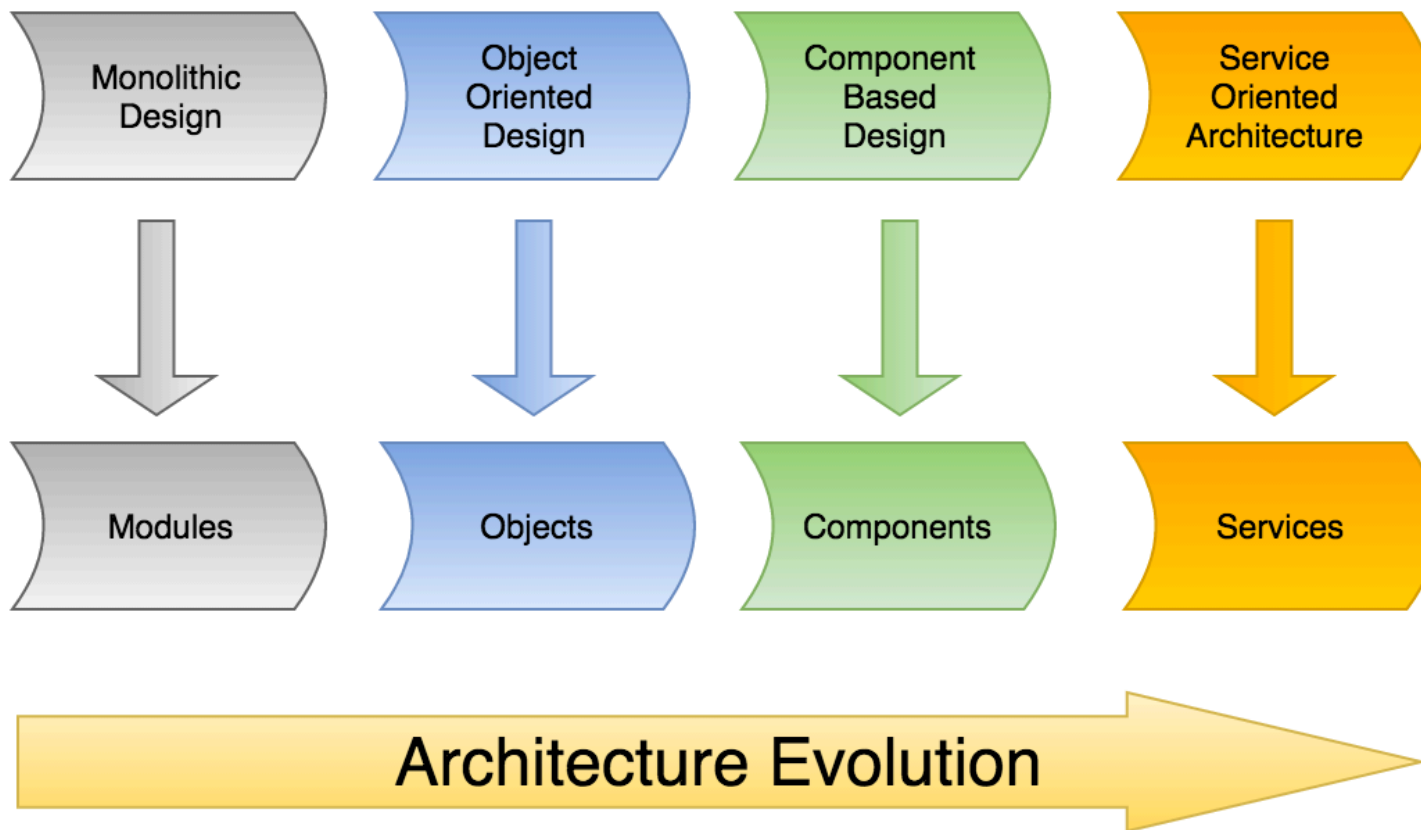
Agenda

- SOA及服务治理介绍
- 服务架构演进
- OCTO架构设计
- 服务治理实践

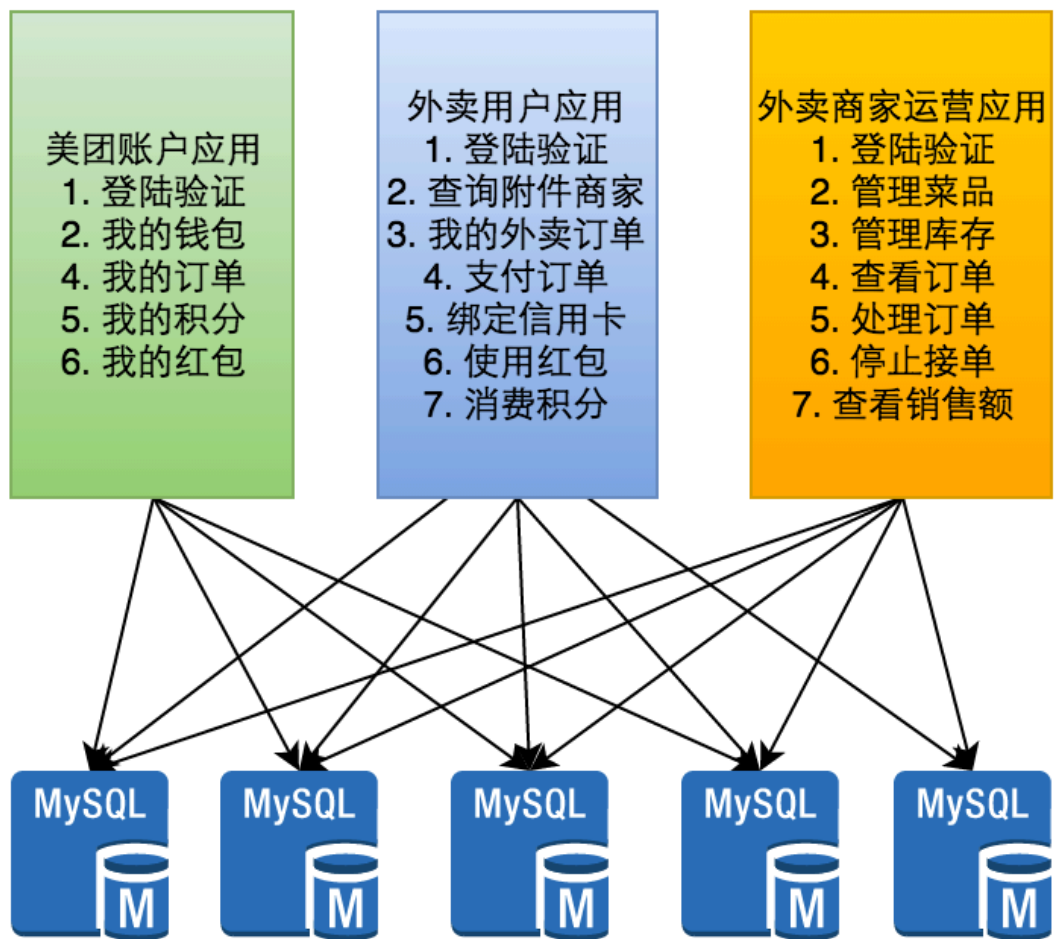
面向服务架构

.....

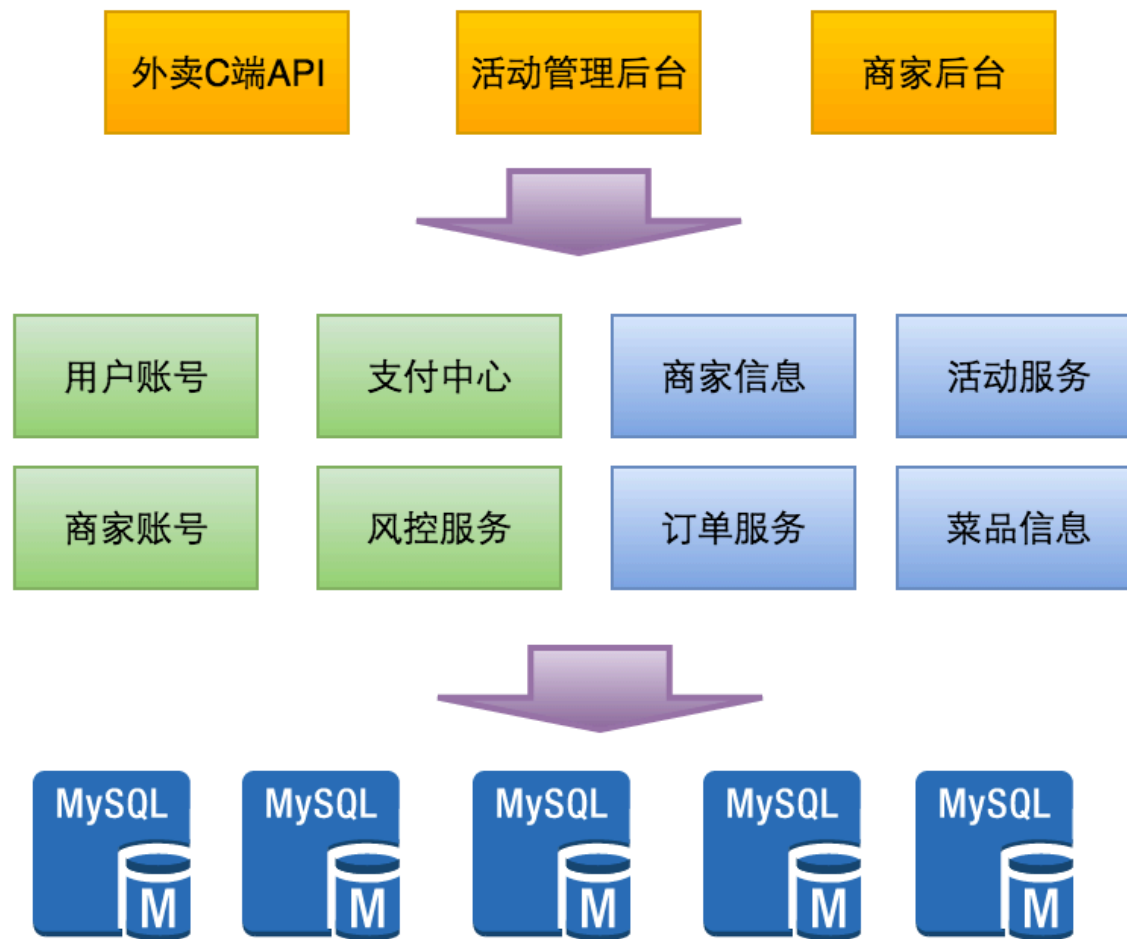
- 单应用设计
 - Function
 - Module
- 面向对象设计
- 组件化设计
- 面向服务架构
 - Service
 - API



Before SOA



After SOA



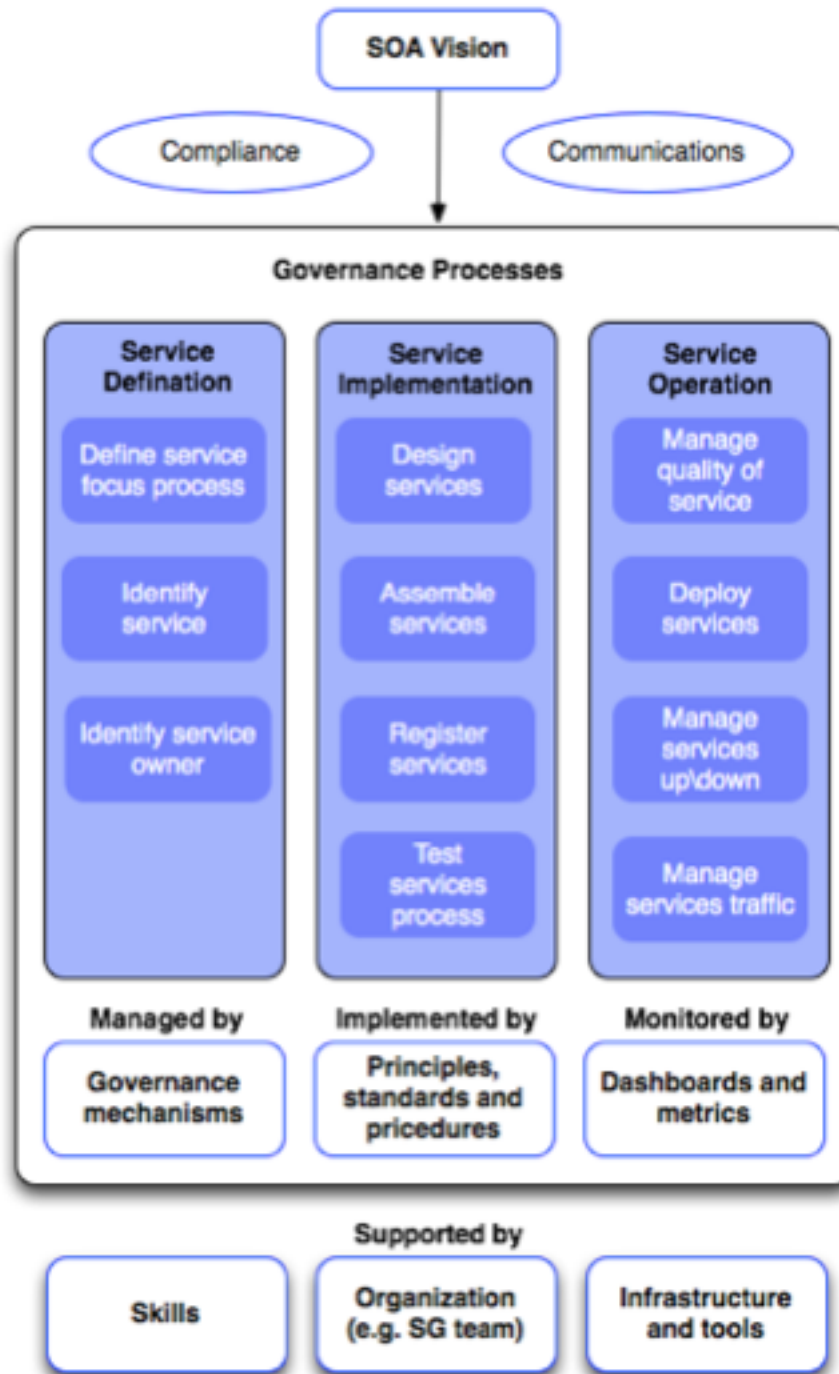
服务治理

- 服务架构原则

- 轻耦合
- 标准化
- 可重用
- 高可靠
- 可伸缩

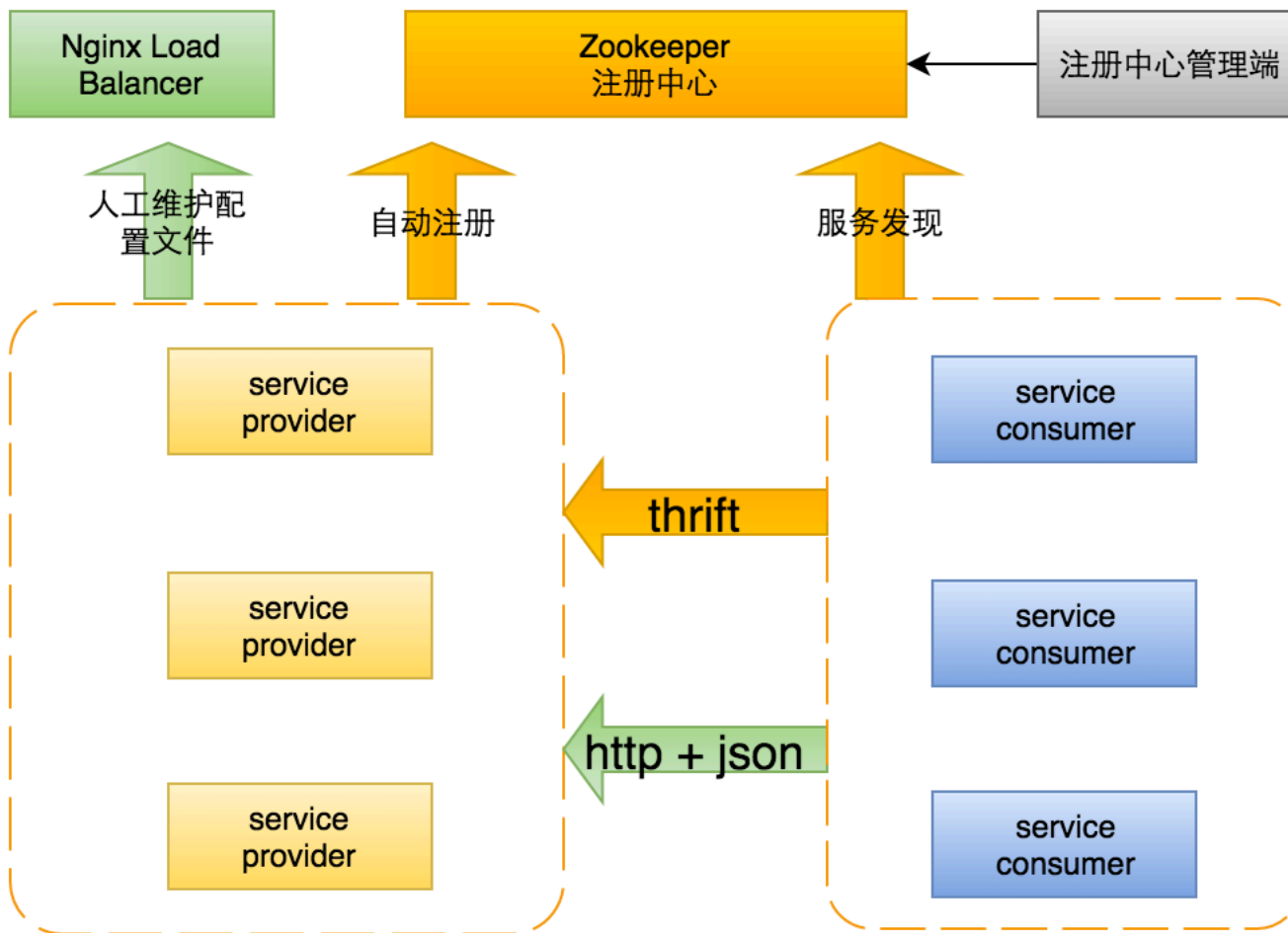
- 服务治理过程

- 服务定义
- 服务实现
- 服务运营



早期服务架构

- 服务注册中心
 - 基于zookeeper
- 服务通信框架
 - Thrift协议
 - 暂只支持Java
- Http协议服务
 - 还需手动维护
- 服务治理涉及不多



早期服务架构 - 遇到的问题

- 服务注册中心
 - 上万服务节点，如何保障足够稳定
 - 多种语言带来的服务注册/发现需求，还存在直接依赖IP的情况
 - 功能扩展的灵活性
- 服务通信框架
 - 过多逻辑放到客户端
 - 大并发高压下，要能保障业务高可用
 - 数据、监控方面支持不够
 - 多语言支持
- 未实现全生命周期的服务治理，做到服务规范化、标准化

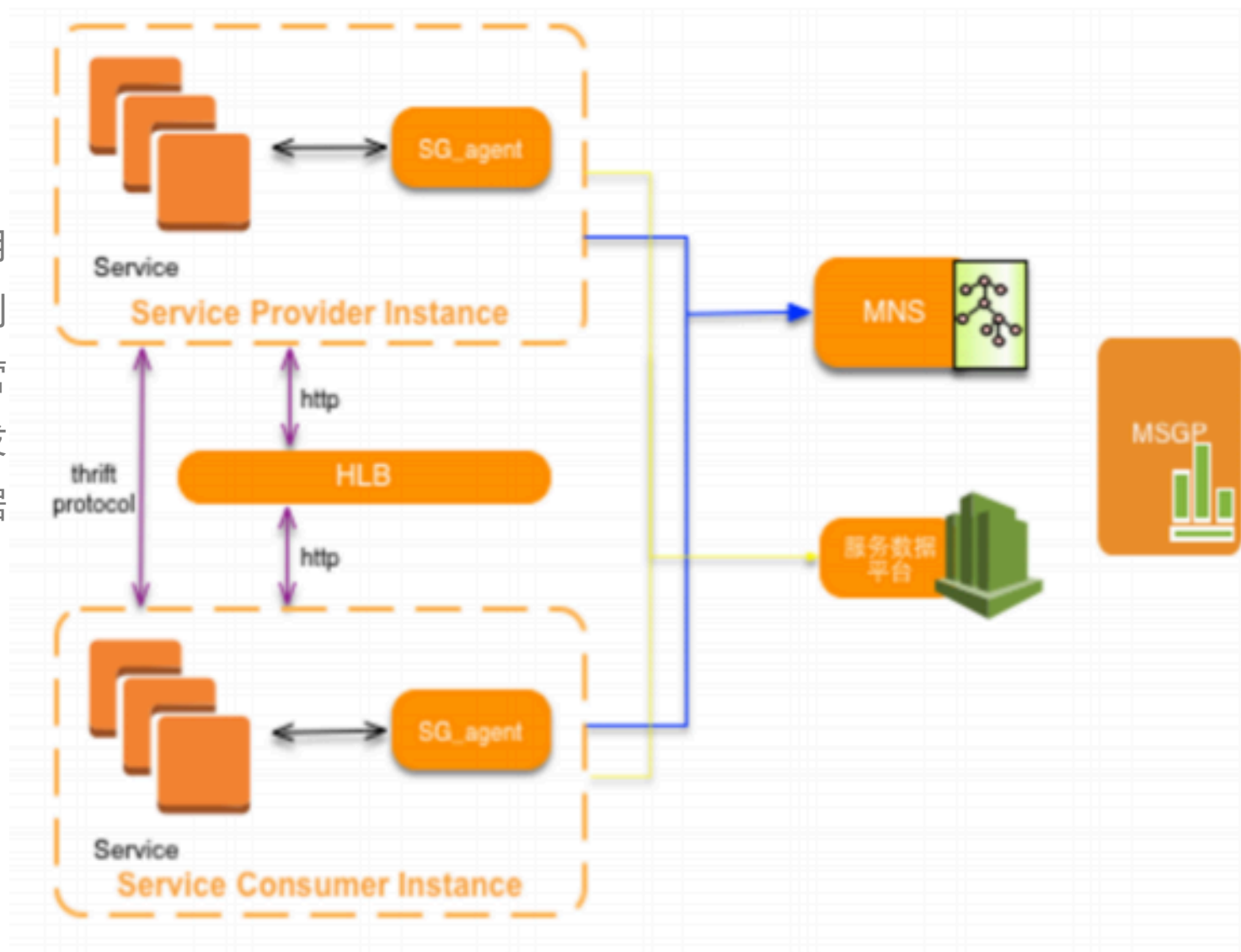
OCTO整体架构

- OCTO是什么？

- 面向服务架构，公司所有业务都使用OCTO作为统一的服务通信框架，达到基本的标准化。并具备良好的服务运营能力，轻松实现服务注册、服务自动发现、负载均衡、容错、灰度发布、数据可视化、监控告警等。

- OCTO关键组件

- MNS：命名服务
- SG_agent：治理代理
- 服务通信框架
- MSGP：服务治理平台



服务注册中心 - 遇到的问题

- zookeeper集群问题

- session timeout , 临时节点抖动下线
- 突发链接过多
- 突发事务过多
- 网络因素导致集群故障
- 其它人为因素 , 如误操作、迁移扩容

可用性百分比	一天24小时为基准	一天8小时为基准
99.9%	8.76小时	2.91小时
99.99%	52.56分钟	17.47分钟
99.999%	5.256分钟	1.747分钟

- 故障恢复及影响

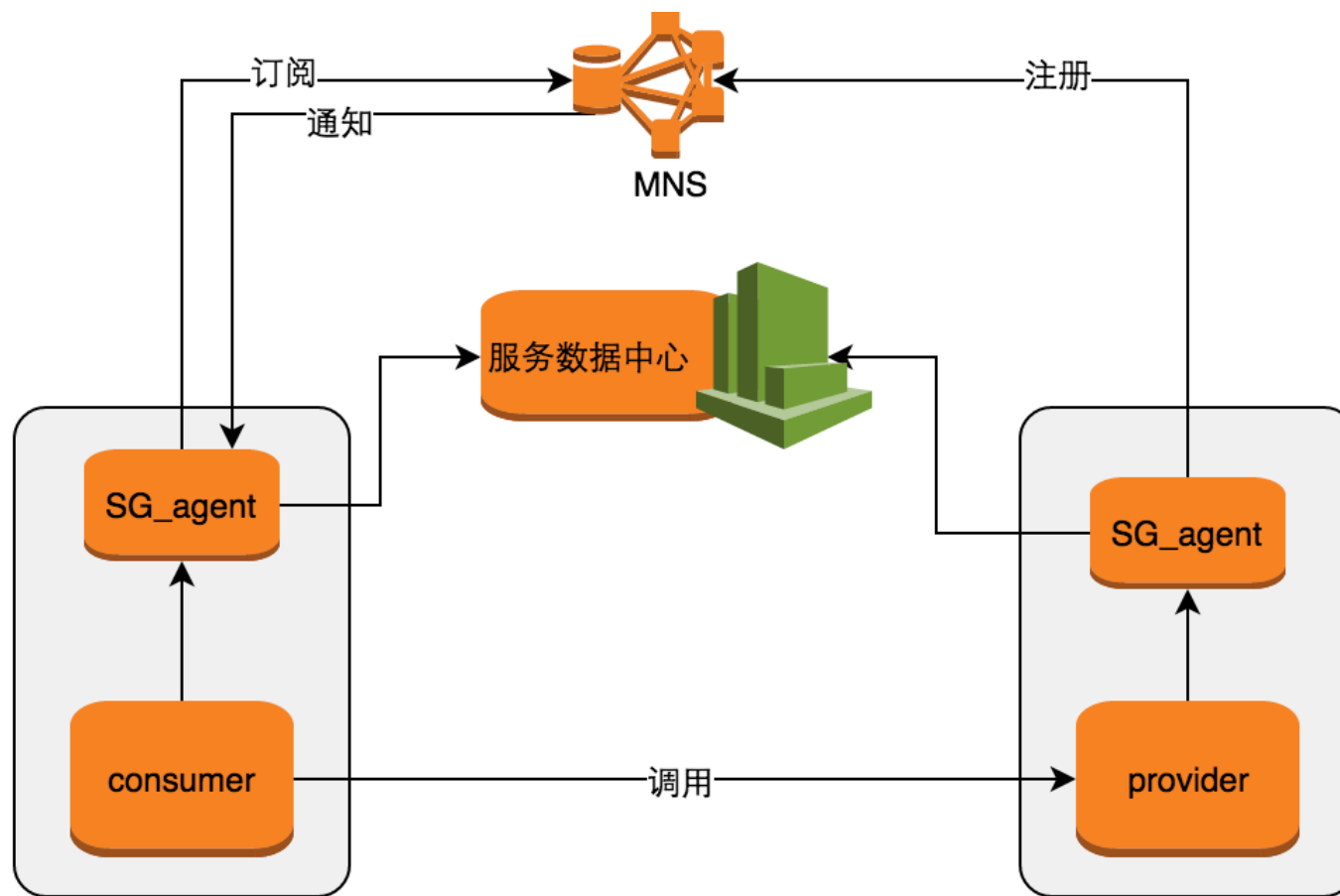
- 集群受业务侧影响导致故障时 , 较难解决
- 网络方面导致的故障 , 恢复时间可能会较长
- 所有业务都需依赖 , 难以隔离影响范围

几种解决思路

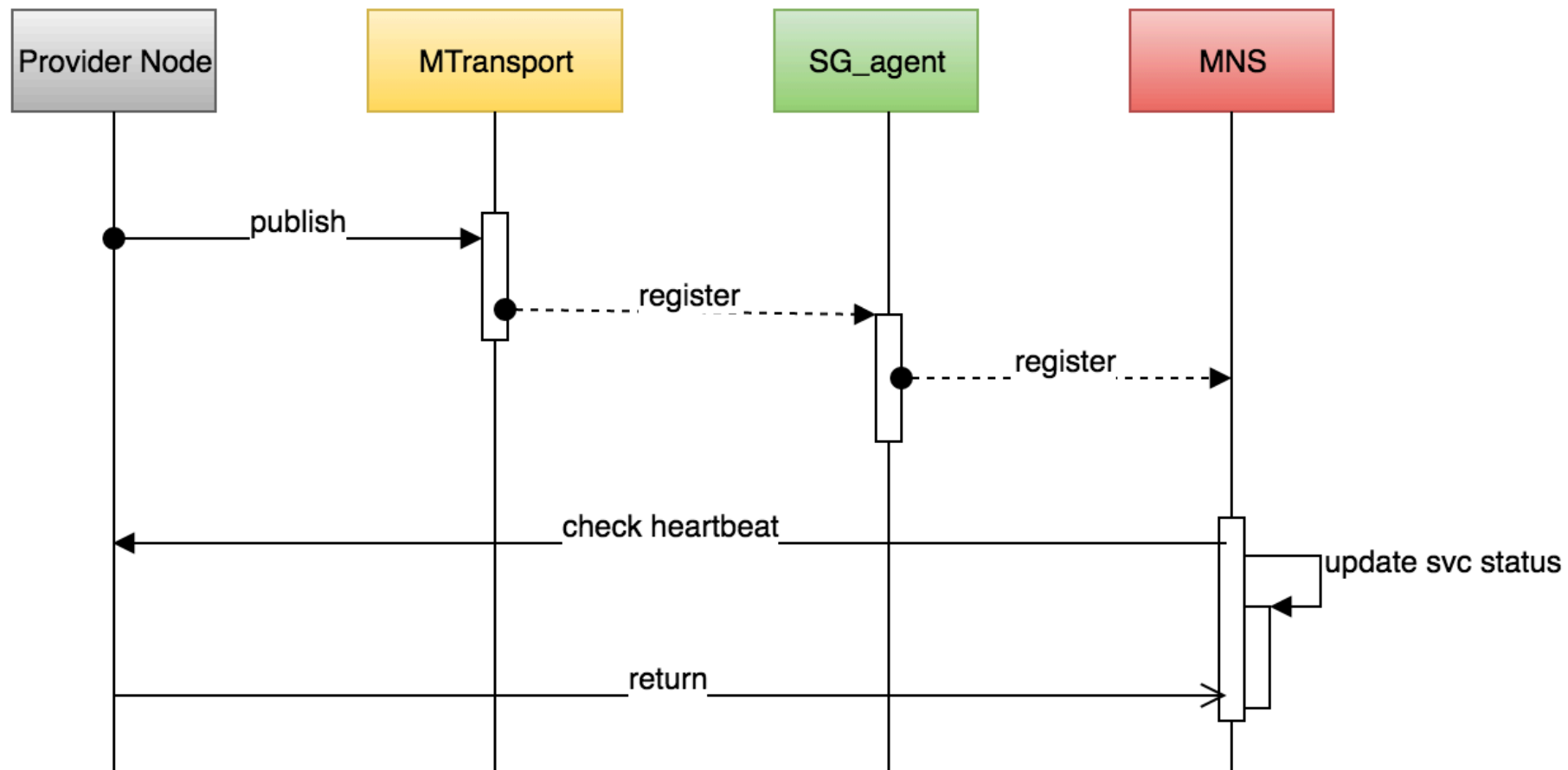
- 增加voter
 - 受限于集群投票机制，通过增加voter并不能很好达到扩容预期
- 增加observer
 - 能很好缓解读请求，但不能解决所有问题
- 集群拆分
 - 能实现事业群间隔离，但带来更多问题
 - 违背统一注册中心的设计初衷
- 框架层内置缓存策略
 - 能部分解决抖动、不可用问题，但服务一旦需要重启还是受影响

OCTO服务注册中心 - 引入SG_agent

- 代理模式
- 本地部署
- 分布式
- 灵活可扩展



引入代理后的注册流程

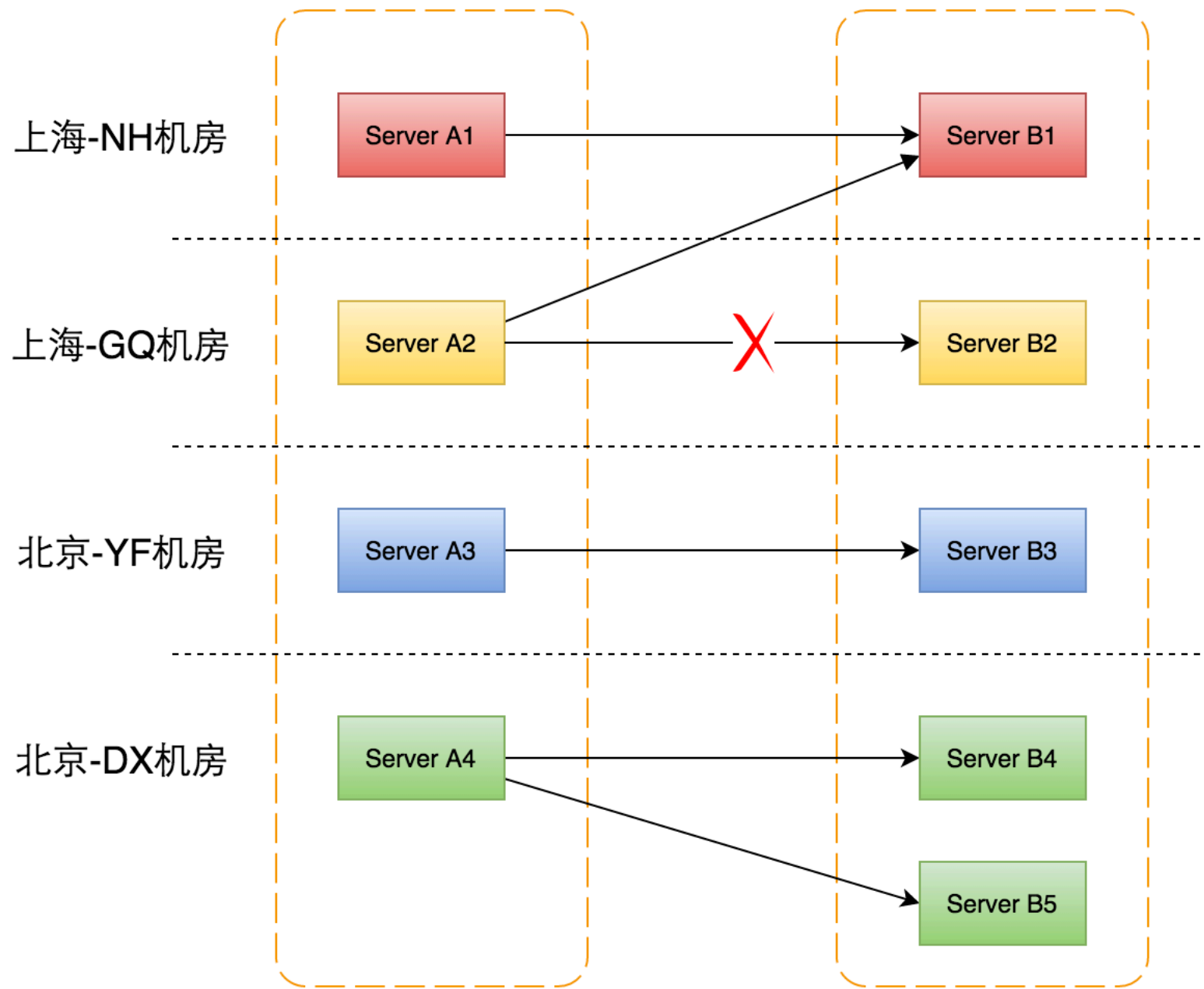


服务通信框架 - 设计要点

- 客户端做薄
- 稳定、高可用
- 数据驱动
- 高性能
- 多语言支持
- 异步化
-

异地多机房下的路由策略需求

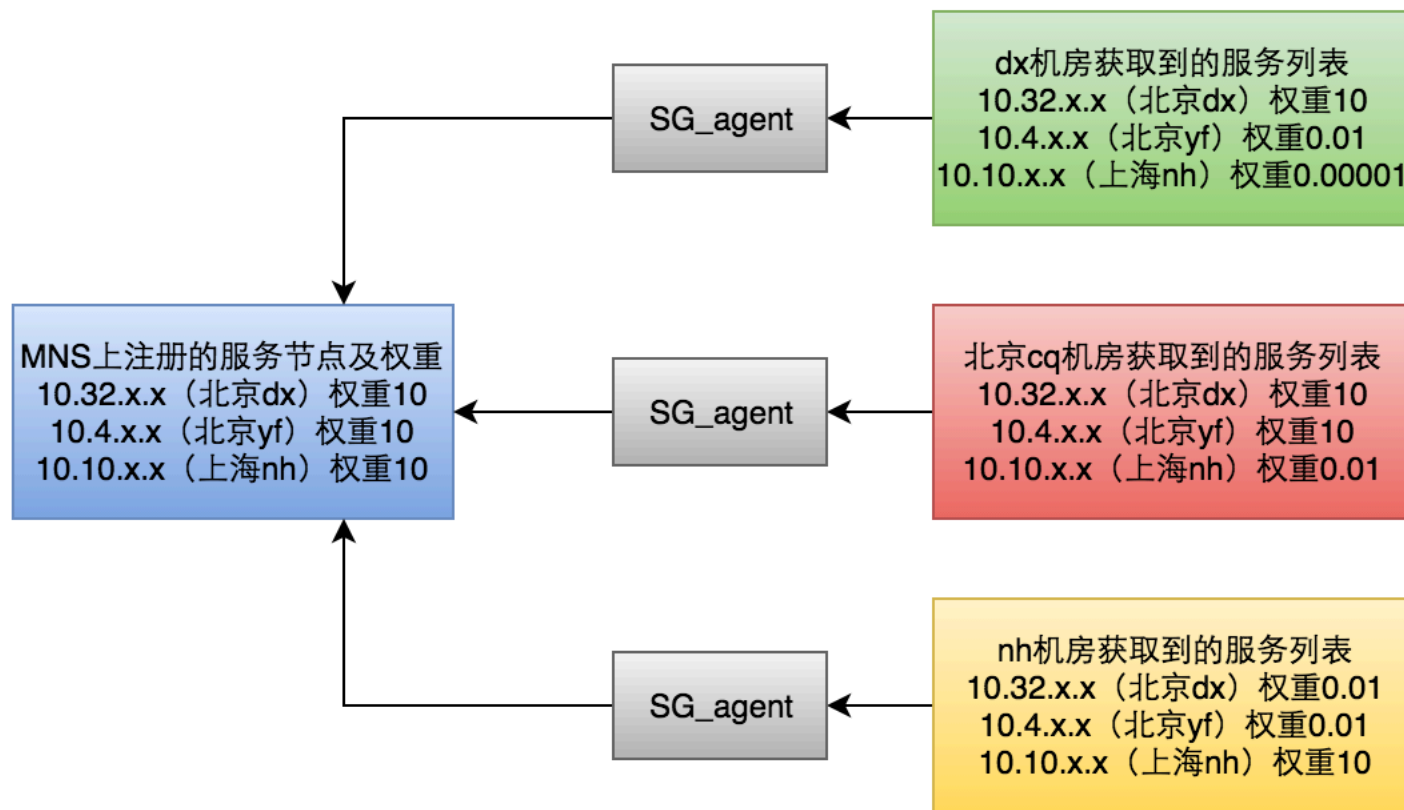
- 单机房场景
 - 全部节点列表，按weight分配流量
- 多机房场景
 - 同机房优先，同机房内的节点，按weight分配流量
 - 同机房节点不可用时，跨机房分配流量
- 异地多机房场景
 - 同机房优先，同机房内的节点，按weight分配流量
 - 同机房节点不可用时，优先在本地域内，跨机房分配流量
 - 本地域无可用节点时，跨地域分配流量



客户端做薄 - 策略下移

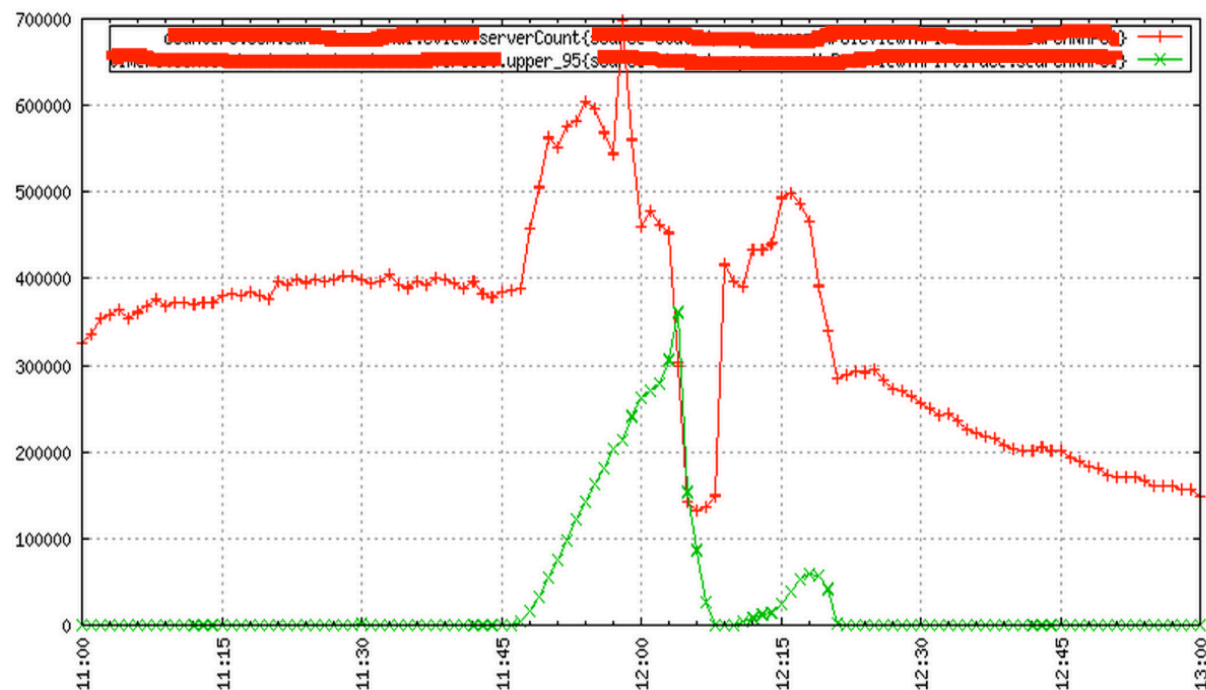
- 框架层实现
 - 单机房 -> 多机房 -> 异地多机房
 - 机房、地域等网络拓扑信息在框架层维护不合适
- 策略下移
 - 权重做归一化处理
 - MNS、SG_agent基于网络拓扑做自动识别
- 下移后的全局拓扑信息可进一步对其它中间件开放
 - KV、MQ等

客户端做薄 - 权重归一化



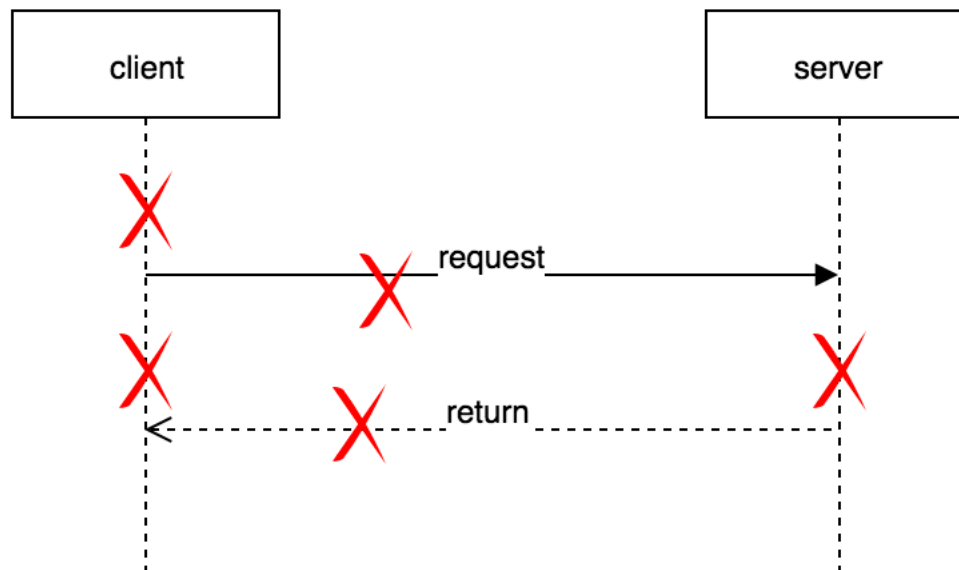
如何保障业务的稳定性、高可用

- 场景1
 - 调用失败了，是否要重试，什么策略？
- 场景2
 - 高峰期出现抖动时，不恰当重试导致雪崩
- 场景3
 - 某个调用方逻辑问题，大量访问，需降级
 - 高峰期压力大，需执行过载保护



导致调用失败的几种场景

- 服务查找阶段failure
 - 无provider，直接抛出
 - 失败节点快速降权、淘汰
- 请求阶段failure
 - 请求确认丢失，可以重试
 - 如果难以确认server是否未收到，需考虑调用是否幂等
- server侧failure
 - 分执行前失败、执行后失败两种情况
 - 从client侧难以区分这两种情况，需考虑调用是否幂等
- client等待response超时
 - 需考虑调用是否幂等



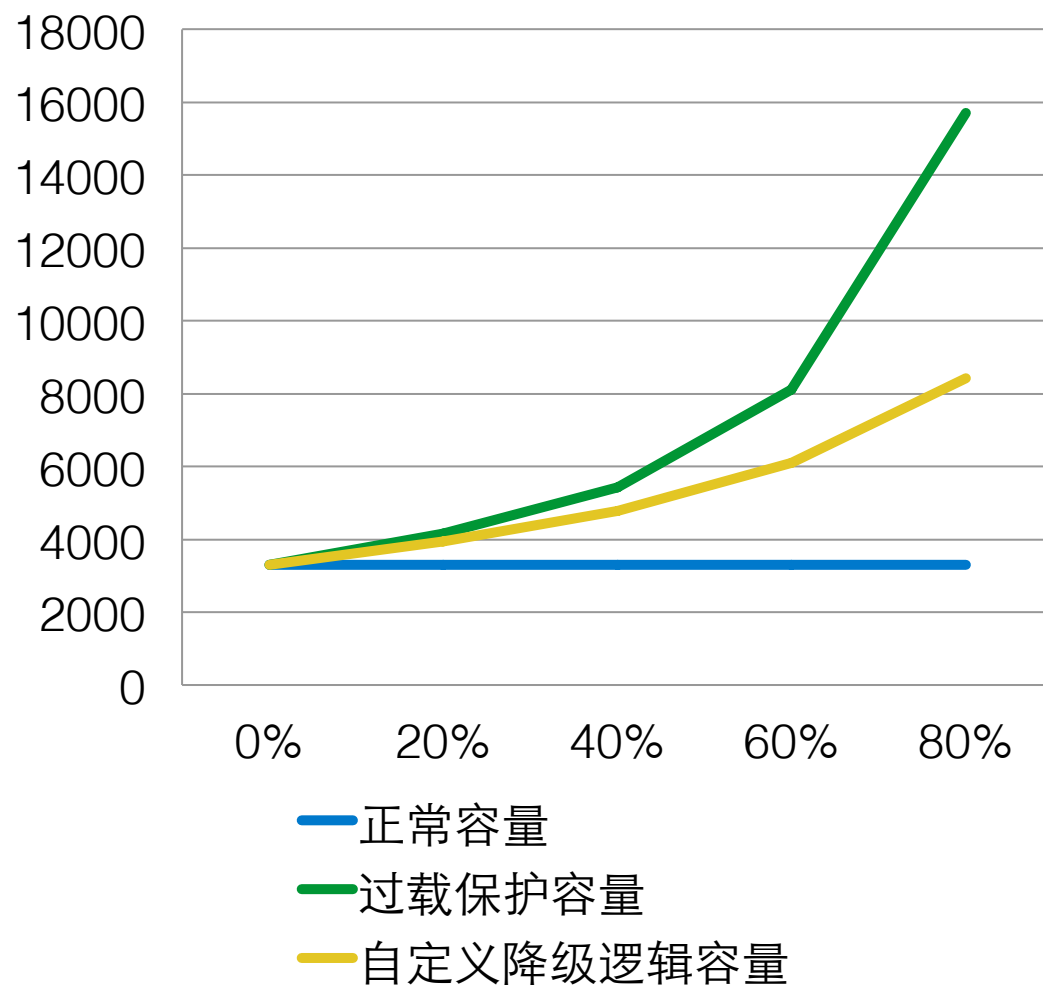
通信框架 - 容错降级机制

- 几种策略

- Fail over : 失败自动重试, 可指定重试次数
- Fail fast : 失败立即报错, 用于非幂等调用

- 过载保护

- 针对来源服务, 一键设置流量的降级比例
- 支持自定义降级逻辑, 不执行一些耗时操作

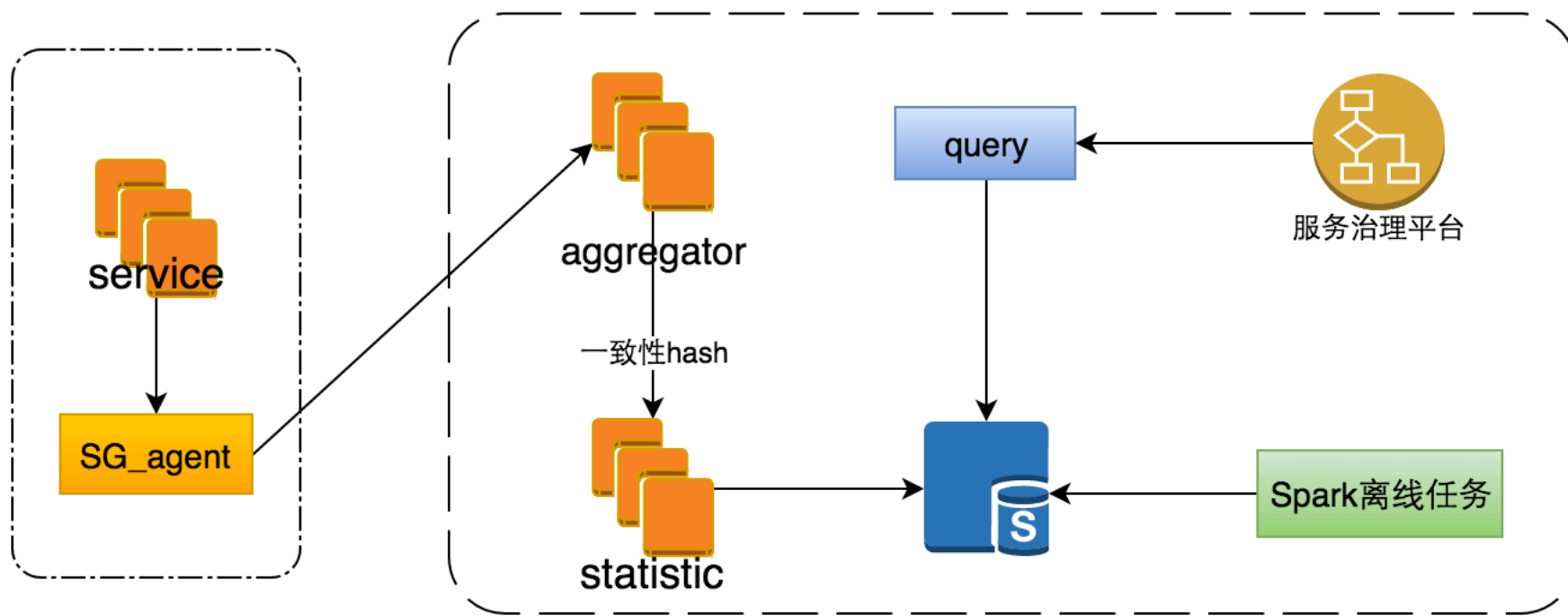


服务治理 - 数据驱动业务

.....

- 谁调用了我？
- 对哪个服务是强依赖？
- 线上调用量突增，来源是哪里？
- 线上响应时间突增，原因是在哪个环节？
- 一次请求失败了，什么原因？怎么能看到详细信息？
- 该业务线涉及几十个服务，整体视图是什么样的？
- 业务线哪些服务是关键服务？可用性如何，当前容量是否足够？

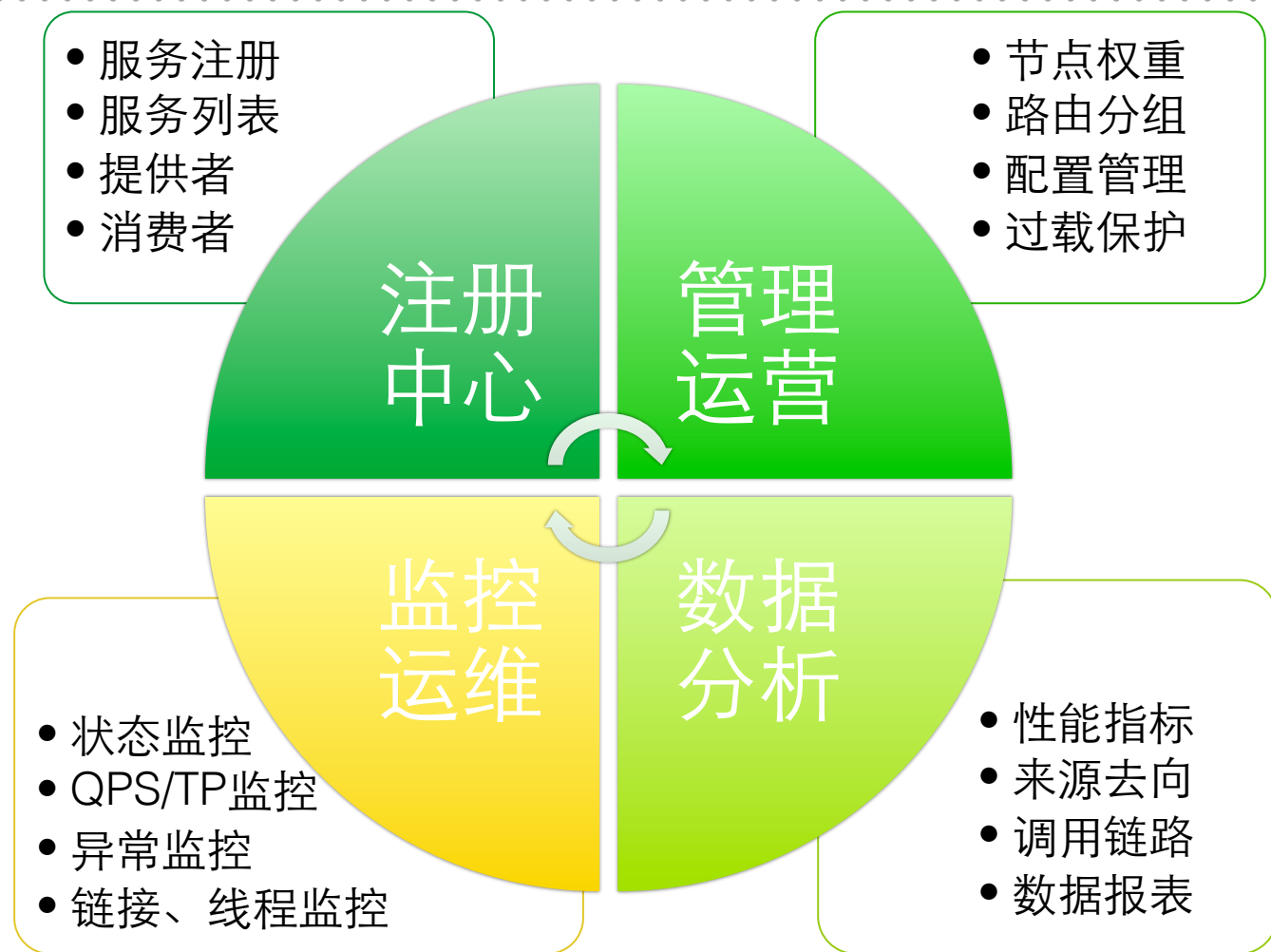
数据采集设计



数据采集&处理

- 自定义协议
 - 框架添加header : requestInfo、responseInfo、traceInfo、context
 - 结合mtrace , 记录调用链路
- 数据采集
 - 采集调用信息后 , 通过SG_agent上报
 - 采集应用层的异常信息 , 实时上报并分析
- 数据处理
 - 一天近百亿的上报数据 , 需要实时处理 : Akka Actor
 - 同一个服务、指标等数据 , 需要先汇聚 : 一致性hash
 - 非实时数据 , 定时分析 : Spark

一站式服务治理平台



服务注册中心

选择服务:

概要

提供者

消费者

实时日志

+新增

类型:

thrift

http

环境?:

prod

stage

test

查询

OCTO服务环境是什么?

主用

备机

启用

禁用

正常

删除

权重

<input type="checkbox"/> 节点	环境	角色	octo版本	权重(点击修改)	状态	工作模式	更新时间	节点状态	操作
<input type="checkbox"/> dx- <div>server02(10.32.106.248):9417</div>	prod	<div>主用</div> <div>备机</div>	mtthrift-v1.6.3	10	正常	octo	2016-06-13 15:56:27	<div>启用</div> <div>禁用</div>	<div>删除</div>
<input type="checkbox"/> dx- <div>server04(10.32.108.249):9417</div>	prod	<div>主用</div> <div>备机</div>	mtthrift-v1.6.3	10	正常	octo	2016-06-13 15:59:29	<div>启用</div> <div>禁用</div>	<div>删除</div>
<input type="checkbox"/> dx- <div>server01(10.32.109.249):9417</div>	prod	<div>主用</div> <div>备机</div>	mtthrift-v1.6.3	10	正常	octo	2016-06-13 15:55:08	<div>启用</div> <div>禁用</div>	<div>删除</div>
<input type="checkbox"/> dx- <div>server03(10.32.109.251):9417</div>	prod	<div>主用</div> <div>备机</div>	mtthrift-v1.6.3	10	正常	octo	2016-06-13 15:57:57	<div>启用</div> <div>禁用</div>	<div>删除</div>

服务数据分析

Dashboard

业务指标

New!

性能指标

上下游分析

New!

来源分析

主机分析

去向分析

调用链分析

New!

角色:

来源

去向

环境:

prod

stage

test

日期:

2016-06-22

查询

可用性指标?

TP耗时数据?

同比环比?

接口	调用总量	成功数/百分比	异常数/百分比	过载数/百分比	QPS(次/秒),环比,同比	TP50耗时(毫秒)	TP90耗时(毫秒),环比,同比	TP95耗时(毫秒)	TP99耗时(毫秒)
all	2447391900	2447391594, 100.0000%	306, 0.0000%	0, 0.0000%	28326.295,-4%,8%	2	4,0%,33%	6	29
www.query.miniserv...	1666041915	1666041719, 100.0000%	196, 0.0000%	0, 0.0000%	19282.893,-2%,8%	2	3,0%,0%	4	9
www.query.miniserv...	294961094	294961062, 100.0000%	32, 0.0000%	0, 0.0000%	3413.902,-2%,9%	2	4,0%,33%	6	13
www.query.miniserv...	177186291	177186264, 100.0000%	27, 0.0000%	0, 0.0000%	2050.767,-13%,11%	2	4,0%,0%	6	11
www.query.miniserv...	54548506	54548502, 100.0000%	4, 0.0000%	0, 0.0000%	631.348,-8%,13%	2	4,0%,0%	6	11



新美大技术团队

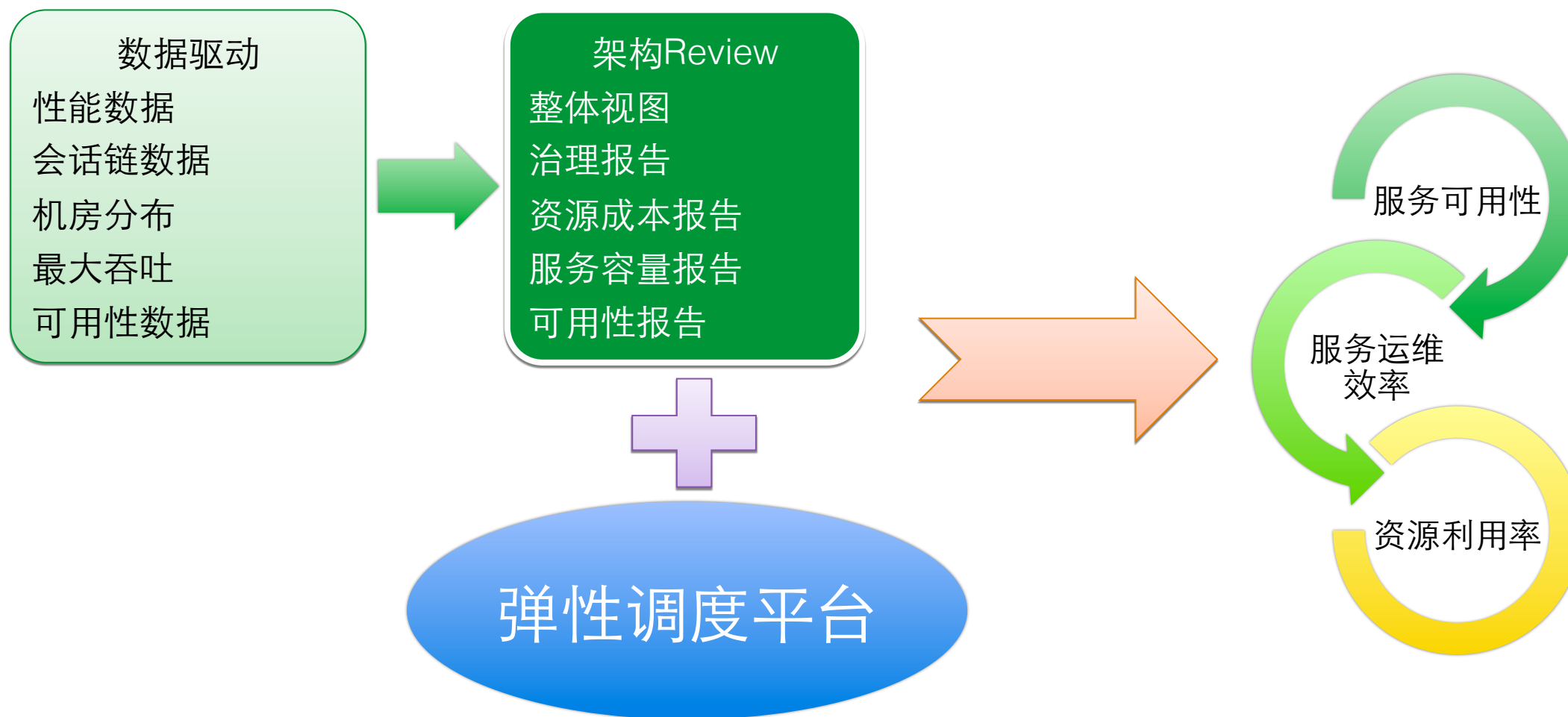
服务视图

配送调度服务 调用 配送人员交易类服务化 的详情

接口	描述	QPS	tp50(ms)	tp90(ms)	tp95(ms)
BmUserThriftIface.getUserByIdAndOtherConditions	该方法暂无描述	1815.433	2	3	3
BmUserThriftIface.isWorkingForRider	该方法暂无描述	58.367	2	3	3
BmOrgThriftIface.getOrgViewById	该方法暂无描述	37.717	2	3	3
BmUserThriftIface.batchGetUserViews	该方法暂无描述	31.317	8	11	11
BmUserThriftIface.searchBmUserOrgSimpleViewPage	该方法暂无描述	18.433	5	8	8
BmUserThriftIface.getVicinalRidersByOrgType	该方法暂无描述	1	13	16	16



OCTO-后续展望



谢谢